

Diagnostics \rightarrow Something wrong in your modeling? p. 7-1

• regression diagnostics: check model assumptions to suggest further improvement after fitting. The building of an empirical model is an iterative process. During the process, it is required to check whether the current fitted model is consistent with data.

• **Q**: what assumptions needed to be checked?



LNp.7-2

model: $Y = X\beta + \epsilon$, $\epsilon \sim N(0, \sigma^2 I)$

- ① \rightarrow error structure: errors independent, equal variance, normally distributed \rightarrow overall pattern
 - ② \rightarrow mean structure: whether $E(Y) = X\beta$ is a correct structure \rightarrow mainly under-fitting or not
 - \rightarrow unusual observations: whether some observations do not fit the model \rightarrow individual pattern
- two types of diagnostic techniques: numerical and graphical

When the linear model is correct

surrogate of ϵ \rightarrow Residual $Y = \hat{Y} + \hat{\epsilon}$

Q: Why cannot get ϵ from $\hat{\epsilon}$? dim(ϵ) = n
dim($\hat{\epsilon}$) = n - p

recall (residuals)

\rightarrow prediction: $\hat{Y} = X(X^T X)^{-1} X^T Y = HY$, H : hat matrix

\rightarrow residuals: $\hat{\epsilon} = Y - \hat{Y} = (I - H)Y = (I - H)X\beta + (I - H)\epsilon = (I - H)\epsilon$

(Note: errors and residuals are different. Q: what difference?)

\rightarrow var($\hat{\epsilon}$) = var($(I - H)\epsilon$) = $(I - H)^2 \sigma^2 = (I - H)\sigma^2 = \sigma^2$

cf. $\text{Var}(\epsilon) = \sigma^2 I$ $\rightarrow E^*[(I - H)\epsilon \epsilon^T (I - H)^T]$

$\sum \hat{\epsilon}_i = 0$ (with intercept)
 $\sum \epsilon_i \neq 0$

$H = [H_{ij}]$

\Rightarrow even though ϵ is uncorrelated and equal variance, $\hat{\epsilon}$ may be not

(Note: H depends on X only) \rightarrow relative variance and correlation of $\hat{\epsilon}$ are known before observing Y in DOE

$\text{Cov}(\hat{\beta}) = (X^T X)^{-1} \sigma^2$
 $\text{Cov}(\hat{Y}) = H \sigma^2$

\rightarrow residuals are useful for detecting lack of fit and checking model assumptions p. 7-2

(Q: Why residuals can do the works?)

examine $\hat{\sigma}^2 \neq \sigma^2$ Q: $\hat{\epsilon}$ for detecting over-fitting? (LNp.6-8)

① $\rightarrow Y = X\beta + \epsilon = \hat{Y} + \hat{\epsilon}$

② true: $Y = X_1\beta_1 + X_2\beta_2 + \epsilon = (X_1\beta_1 + H_1 X_2\beta_2) + ((I - H_1)X_2\beta_2 + \epsilon) = \hat{Y}_{X_1} + \hat{\epsilon}_{X_1}^*$

lacked of fit \rightarrow fitted: $Y = X_1\beta_1 + \epsilon^*$

Leverage $\rightarrow X_1 [\beta_1 + (x_i^T x_i)^{-1} x_i^T X_2 \beta_2]$ check LNp.5-9

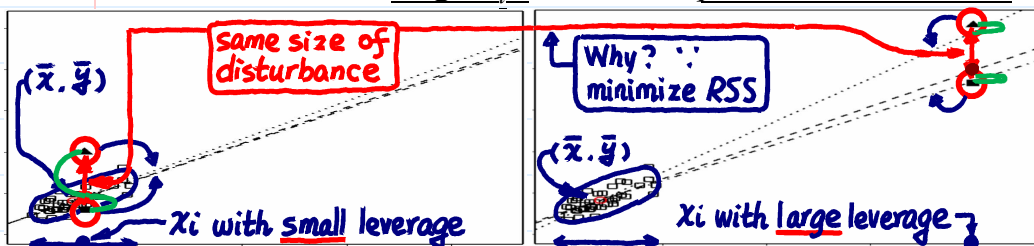
• leverage: $h_i \equiv H_{ii}$ (Note 1. var($\hat{\epsilon}_i$) = $(1 - h_i)\sigma^2$. Note 2. h_i is known before observing Y)

(x_i, y_i)
 $\leftrightarrow \hat{\epsilon}_i$
 $\leftrightarrow h_i$

x_i whose h_i is large \Rightarrow var($\hat{\epsilon}_i$) small \Rightarrow fitted model has to force to fit close to y_i

x_i whose h_i is small \Rightarrow var($\hat{\epsilon}_i$) large \Rightarrow in this x_i , model cannot fit so well

- \rightarrow h_i roughly determines how close (x_i, y_i) to the regression surface (i.e., (x_i, \hat{y}_i))
- \rightarrow observations with large h_i 's should be paid more attention. (Q: why?) Q: y_i not so close to \hat{y}_i



Q: why x_i with large leverage has stronger influence on fit?

\rightarrow for linear model with an intercept, its fitted model must pass the point (\bar{x}, \bar{y})

$E(y) = \beta_0 + \sum_{i=1}^{p-1} \beta_i x_i = \left(\beta_0 + \sum_{i=1}^{p-1} \beta_i \bar{x}_i \right) + \sum_{i=1}^{p-1} \beta_i (x_i - \bar{x}_i) = \beta'_0 + \sum_{i=1}^{p-1} \beta_i (x_i - \bar{x}_i)$

$\Rightarrow \hat{\beta}'_0 = \bar{y} \Rightarrow \hat{E}(y) = \bar{y} + \sum_{i=1}^{p-1} \hat{\beta}_i (x_i - \bar{x}_i)$

orthogonality also, check LNp.3-7 centering each X 力矩 = 力臂 * 作用力

Q: why is it called leverage?